

BEST AVAILABLE COPY

REMARKS

Claims 42-50, 67, 69-72, 74, 75, 80, 81, 88, 90-93, and 95-104 are pending in the application. Claims 42-44, 46-48, 67, 69, 70, 72, 74, 75, 80, 81, 88, 90-93, and 95-99 have been amended. Claims 42-45, 67, 69, 70, 80, 88, 90-93, and 95-104 have been withdrawn from consideration. Claims 1-41, 51-66, 68, 73, 76-79, 82-87, 89, and 94 have been cancelled without prejudice. These amendments add no new matter.

Restriction/Election

In response to the Restriction Requirement, Applicants elect the invention of Group III, drawn to nucleic acids encoding the LBP-2 polypeptide of SEQ ID NO:43 and variants and fragments thereof. The election is made with traverse.

The Examiner divided the claims into nine separate groups. For the reasons provided below, applicants respectfully request that the claims of Groups II, VI, VII, and IX (partially) be examined together with the claims of Group III in the present application.

The pending claims are directed to nucleic acids encoding a novel LDL-binding polypeptide ("LBP") termed LBP-2. Both human and rabbit LBP-2 sequences are recited in the claims. Groups II, III, VI, VII, and IX (partially) are directed to human LBP-2 sequences. Groups I, IV, V, VIII, and IX (partially) are directed to non-elected rabbit LBP-2 sequences.

The human LBP-2 polypeptide is described in SEQ ID NO:43 (full length polypeptide) and SEQ ID NO:7 (amino acids 322-538 of human LBP-2). The human LBP-2 polypeptide of SEQ ID No:7 (Group II) is a fragment of the elected SEQ ID NO:43. SEQ ID NO:16 (Group VI) is a specific nucleotide sequence encoding SEQ ID NO:7, and SEQ ID NO:45 (Group VII) is a specific nucleotide sequence encoding SEQ ID NO:43. SEQ ID NOS 30, 31, 32, and 33 (Group IX) each encodes a polypeptide fragment of the elected SEQ ID NO:43.

Because of the extremely high sequence relatedness between SEQ ID NO:43 and SEQ ID NO:7, applicants submit that prosecution will be facilitated by the simultaneous examination of nucleic acids encoding each of these human LBP-2 polypeptides. In addition, the issues raised during the course of prosecution of these human LBP-2 nucleic acid sequences are expected to

be similar and, therefore, simultaneous examination is not expected to be unduly burdensome. Applicants also note that in a parent application of the present application (U.S. Patent No. 6,632,923), all of the human LBP-2 polypeptides were examined together with the respective variants and fragments. Applicants have cancelled claims directed to rabbit LBP-2 nucleic acids from the present application.

In light of the above comments, applicants respectfully request that the Examiner examine the human LBP-2 nucleic acid sequences of Groups II, III, VI, VII, and IX.

35 U.S.C. § 112, 1st Paragraph (Enablement)

On pages 8-9 of the Office Action, the Examiner rejected claims 59-61, 72, 74, and 75 as allegedly "containing subject matter which was not described in the specification in such a way as to enable one skilled in the art to which it pertains, or with which it is most nearly connected, to make and/or use the invention." In particular, the Examiner stated that

the specification fails to describe or provide guidance about the nucleotide sequence that encodes a polypeptide comprising an amino acid sequence having identity to a fragment of at least 10 or 20 or 30 amino acid residues of the encoded polypeptide of SEQ ID NO: 43 (claims 72, 74, 75). It is not clear to a skilled artisan that what is the position of these 10, 20, and 30 amino acids in relation to the amino acid sequence set forth in SEQ ID NO: 43. Although Examples 2, 3, 4, 5 (pages 40-45) demonstrate the full-length cDNA encoding LDL binding protein, this is not demonstrative of any fragments or analogs that are claimed in claims 59-61 and claims 72, 74, and 75. For these reasons it would require undue experimentation to make the claimed invention.

Claims 59-61 have been cancelled thereby rendering their objection moot.

Applicants traverse the rejection of claims 72, 74, and 75 in light of the claim amendments and the following comments.

As amended, claims 72, 74, and 75 are directed to an isolated nucleic acid comprising a nucleotide sequence that encodes a polypeptide comprising an amino acid sequence that binds to LDL and is identical to a fragment of at least 10, 20, or 30 amino acid residues of the human LBP-2 polypeptide of the SEQ ID NO:43. A person skilled in the biological arts knows how to generate nucleic acids encoding fragments of the LBP-2 polypeptide and how to test the ability of such polypeptide fragments to bind to LDL. For example, as detailed in the specification

(see, e.g., page 21), fragments of a polypeptide can be generated by removing one or more nucleotides from one end (for a terminal fragment) or both ends (for an internal fragment) of a nucleic acid that encodes a polypeptide. Nucleic acids that encode fragments of a polypeptide can also be generated by, e.g., random shearing, endonuclease restriction digestion, or a combination of any of these methods. Expression of such a recombinant DNA would produce the desired LBP-2 fragments.

The specification instructs how to evaluate the ability of LBP-2 polypeptide fragments to bind to LDL by using methods such as affinity chromatography, affinity coelectrophoresis, or ELISA (see specification at page 21, line 2 to page 22, line 3). Consistent with the preceding comments, Example 8 indicates that a particular stretch of acidic amino acids of the human LBP-2 (about amino acids 329-354) participates in the binding of LBP-2 to LDL (page 49, lines 4-20). Examples 9 and 10 further detail methods for determining whether LBP-2 polypeptides bind to LDL in the presence of a given candidate inhibitor (page 49, line 22 to page 51, line 10).

In light of the foregoing, a person of ordinary skill in the biological arts would have been able to make and use an isolated nucleic acid that encodes an LDL-binding fragment of LBP-2 without undue experimentation and with a reasonable expectation of success.

At pages 9-12 of the Office Action, the Examiner rejected claims 46-48, 59-61, 72, 73-75, 81, and 85-87 as allegedly not enabled. In particular, the Examiner stated that

the specification, while being enabling for an isolated nucleic acid comprising a sequence that encodes a polypeptide of an amino acid sequence set forth in SEQ ID NO:43 that binds to low density lipoprotein (LDL); does not reasonably provide enablement for all the LDL binding proteins, and fragments and mutants generated from any position located on the sequence of SEQ ID NO:43. The specification does not enable any person skilled in the art to which it pertains, or with which it is most nearly connected, to make and/or use the invention commensurate in scope with these claims. The specification, however, only discloses cursory conclusions (see page 8-24) to support the findings.

Claims 59-61, 73, and 85-87 have been cancelled thereby rendering their rejection moot.

Applicants respectfully traverse the rejection of claims 46-48, 72, 74-75, and 81 in light of the claim amendments and the comments provided below.

The claims rejected herein are directed to variants or fragments of human LBP-2 that retain the ability to bind LDL. It is well within the grasp of the biologist of ordinary skill to prepare, for example, a polypeptide having at least 80%, at least 90%, at least 95%, or at least 98% sequence identity to the human LBP-2 of SEQ ID NO:43. The specification details standard mutagenesis methods that can be used to make amino acid sequence variants (page 19, line 2 to page 20, line 2). Furthermore, the specification instructs, and the skilled biologist is well aware, that conservative amino acid substitutions can be made in the LBP-2 polypeptide sequence so as to reduce the likelihood that a given amino acid sequence will result in a loss of LBP-2 function (page 17, line 21 to page 18, line 15). In addition, fragments of the full-length LBP-2 polypeptide can be generated by using standard techniques that are detailed above and in the specification (e.g., page 21, line 2 to page 22, line 3).

In addition to being able to readily produce nucleic acids encoding human LBP-2 fragments or sequence variants, it would have required no undue experimentation for the skilled artisan to identify those variants that retain the specific LDL binding activity recited in the claims 46-48, 74-75, and 81. By using the assays described in the specification (e.g., page 47, line 23 to page 49, line 20), the skilled artisan would have been able to determine, without undue experimentation and with a reasonable expectation of success, whether a given human LBP-2 fragment or sequence variant binds to LDL.

At page 11 of the Office Action, the Examiner stated that the "specification has provided no guidance to enable one of ordinary skill in the art to determine, without undue experimentation, the positions in the protein, which are tolerant to change (e.g., by amino acid deletions, insertions, or substitutions) and the nature and extent of changes that can be made in these positions."

Although it is possible in certain cases to abolish the functional activity of a protein by mutating a critical amino acid residue, this does not mean that one of ordinary skill cannot nonetheless readily make functional analogs of a given protein (e.g., LBP-2) without undue

experimentation. In fact, as detailed in the enclosed publication of Bowie et al. (1990) Science 247:1306-10 ("Exhibit A"), "proteins are surprisingly tolerant of amino acid substitutions." Exhibit A cites as evidence of this assertion a study carried out on the *lac* repressor that found that of approximately 1500 single amino acid substitutions at 142 positions in the protein, "about one-half of all substitutions were phenotypically silent." Thus, one can expect, based on Exhibit A's disclosure, that a significant percentage of random substitutions in a given protein will result in mutated proteins with full or nearly full activity. These are far better odds than those at issue in *In re Wands*, 858 F.2d 731 (Fed. Cir. 1988), cited by the Examiner on page 9, in which the court found that screening many hybridomas to find the few that fell within the claims was not undue experimentation. The question is not whether it is possible to abolish activity of a given protein by introducing a point mutation, but rather whether one of ordinary skill can produce, without undue experimentation, mutants in which the activity is not abolished.

Based on Exhibit A's disclosure, one would predict that even random substitution of amino acid residues of a human LBP-2 polypeptide would result in a large pool of mutants having the LDL binding activity recited in the claims. Furthermore, as detailed herein, the specification amply teaches the skilled artisan how to select those mutants having the activity required by the claims. In light of these comments, Applicants submit that one of ordinary skill in the art would have been able, at the filing of the present application, to make and use the claimed nucleic acids without undue experimentation. Accordingly, Applicants request that the Examiner withdraw the rejection.

35 U.S.C § 112, 2nd Paragraph (Indefiniteness)

Claims 46-48 and 59-61 were rejected as allegedly indefinite in their use of the term "identical." Claims 59-61 have been cancelled thereby rendering their rejection moot. For claims 46-48, Applicants have adopted the Examiner's suggested language by directing the claims to an isolated nucleic acid comprising a nucleotide sequence that encodes a polypeptide comprising an amino acid sequence and has at least 80%, 90%, or 95% "sequence identity" to

the sequence of SEQ ID No:43. In light of these amendments, Applicants request that the Examiner withdraw the rejection.

The Examiner rejected independent claims 46, 59, 73, 81, and 85 as allegedly indefinite because of the use of the term "LDL" in the absence of the fully spelled out phrase "low density lipoprotein." Claim 42 has been amended to provide the fully spelled out term, which is followed by the acronym throughout the remainder of the claims. In light of these amendments, Applicants request that the Examiner withdraw the rejection.

35 U.S.C. 102(e) (Anticipation)

The Examiner rejected claims 72, 73, and 85 as allegedly anticipated by Colasanti et al., U.S. Patent No. 6,177,614 ("Colasanti"). According to the Examiner,

Colasanti's peptide is considered for the encoded peptide sequence fragment of at least 10 amino acid residues of SEQ ID NO:43 (claims 72, 85). Colasanti's peptide having the structure of the claimed encoded peptide of instant application considered anticipating the LDL binding of the claimed peptide (73). Therefore, claims 72, 73, and 85 of the instant application are being anticipated by Colasanti et al.

Claims 73 and 85 have been cancelled, rendering the objection moot for these claims.

Applicants traverse the rejection of claim 72 in light of the claim amendments and the following comments.

Claim 72, as amended, is directed to an isolated nucleic acid comprising a nucleotide sequence that encodes a polypeptide comprising an amino acid sequence that binds to LDL and is identical to a fragment of at least ten amino acids of SEQ ID NO:43.

Colasanti discloses the Id gene in maize plants, which is similar to that of genes encoding zinc-finger regulatory proteins in animals. However, nothing in Colasanti suggests that the Id gene encodes a protein that binds LDL.

"The fact that a certain result or characteristic may occur or be present in the prior art is not sufficient to establish the inherency of that result or characteristic" (MPEP § 2112, citing In re Rijckaert, 9 F.3d 1531, 1534 (Fed. Cir. 1993)) (emphasis in original). To rely on inherency, "the examiner must provide a basis in fact and/or technical reasoning to reasonably support the

determination that the allegedly inherent characteristic necessarily flows from the teachings of the applied prior art.” (MPEP § 2112, citing Ex parte Levy, 17 USPQ2d 1461, 1464 (Bd. Pat. App. & Inter. 1990)) (emphasis in original). “Inherency, however, may not be established by probabilities or possibilities.” (MPEP § 2112, citing In re Robertson, 169 F.3d 743, 745 (Fed. Cir. 1999)).

There is no evidence of record that would lead the skilled artisan to conclude that a peptide disclosed by Colasanti binds to LDL. Because there is no basis in fact or technical reasoning to reasonably conclude that Colasanti discloses an isolated nucleic acid sequence that encodes a polypeptide comprising an amino acid sequence that is identical to a fragment of at least 10 amino acids of SEQ ID No:43 and also binds to LDL, the reference does not anticipate claim 72. Applicants request that the Examiner withdraw the rejection.

CONCLUSIONS

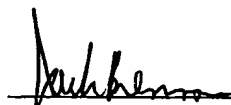
Applicants submit that all grounds for rejection have been overcome, and that all claims are now in condition for allowance.

Enclosed is a Petition for Two Month Extension of Time and a check for the Petition for Extension of Time fee. Please apply any other charges or credits to deposit account 06-1050, referencing Attorney Docket No. 10797-004002.

Respectfully submitted,

Date:

September 9, 2004



Jack Brennan
Reg. No. 47,443

Fish & Richardson P.C.
45 Rockefeller Plaza, Suite 2800
New York, New York 10111
Telephone: (212) 765-5070
Facsimile: (212) 258-2291

Deciphering the Message in Protein Sequences: Tolerance to Amino Acid Substitutions

JAMES U. BOWIE,* JOHN F. REIDHAAR-OLSON, WENDELL A. LIM,
ROBERT T. SAUER

amino acid sequence encodes a message that determines the shape and function of a protein. This message is highly degenerate in that many different sequences can be for proteins with essentially the same structure and activity. Comparison of different sequences with similar messages can reveal key features of the code and improve understanding of how a protein folds and how it performs its function.

THE GENOME IS MANIFEST LARGELY IN THE SET OF PROTEINS that it encodes. It is the ability of these proteins to fold into unique three-dimensional structures that allows them to function and carry out the instructions of the genome. Thus, comprehending the rules that relate amino acid sequence to structure is fundamental to an understanding of biological processes. Because an amino acid sequence contains all of the information necessary to determine the structure of a protein (1), it should be possible to predict structure from sequence, and subsequently to infer detailed aspects of function from the structure. However, both problems are extremely complex, and it seems unlikely that either can be solved in an exact manner in the near future. It may be possible to obtain approximate solutions by using experimental data to simplify the problem. In this article, we describe how an analysis of allowed amino acid substitutions in proteins can be used to reduce the complexity of sequences and reveal important aspects of structure and function.

Methods for Studying Tolerance to Sequence Variation

There are two main approaches to studying the tolerance of an amino acid sequence to change. The first method relies on the process of evolution, in which mutations are either accepted or rejected by natural selection. This method has been extremely successful for proteins such as the globins or cytochromes, for which sequences from many different species are known (2-7). The second approach uses genetic methods to introduce amino acid changes at

specific positions in a cloned gene and uses selections or screens to identify functional sequences. This approach has been used to great advantage for proteins that can be expressed in bacteria or yeast, where the appropriate genetic manipulations are possible (3, 8-11). The end results of both methods are lists of active sequences that can be compared and analyzed to identify sequence features that are essential for folding or function. If a particular property of a side chain, such as charge or size, is important at a given position, only side chains that have the required property will be allowed. Conversely, if the chemical identity of the side chain is unimportant, then many different substitutions will be permitted.

Studies in which these methods were used have revealed that proteins are surprisingly tolerant of amino acid substitutions (2-4, 11). For example, in studying the effects of approximately 1500 single amino acid substitutions at 142 positions in *lac* repressor, Miller and co-workers found that about one-half of all substitutions were phenotypically silent (11). At some positions, many different, nonconservative substitutions were allowed. Such residue positions play little or no role in structure and function. At other positions, no substitutions or only conservative substitutions were allowed. These residues are the most important for *lac* repressor activity.

What roles do invariant and conserved side chains play in proteins? Residues that are directly involved in protein functions such as binding or catalysis will certainly be among the most conserved. For example, replacing the Asp in the catalytic triad of trypsin with Asn results in a 10^4 -fold reduction in activity (12). A similar loss of activity occurs in λ repressor when a DNA binding residue is changed from Asn to Asp (13). To carry out their function, however, these catalytic residues and binding residues must be precisely oriented in three dimensions. Consequently, mutations in residues that are required for structure formation or stability can also have dramatic effects on activity (10, 14-16). Hence, many of the residues that are conserved in sets of related sequences play structural roles.

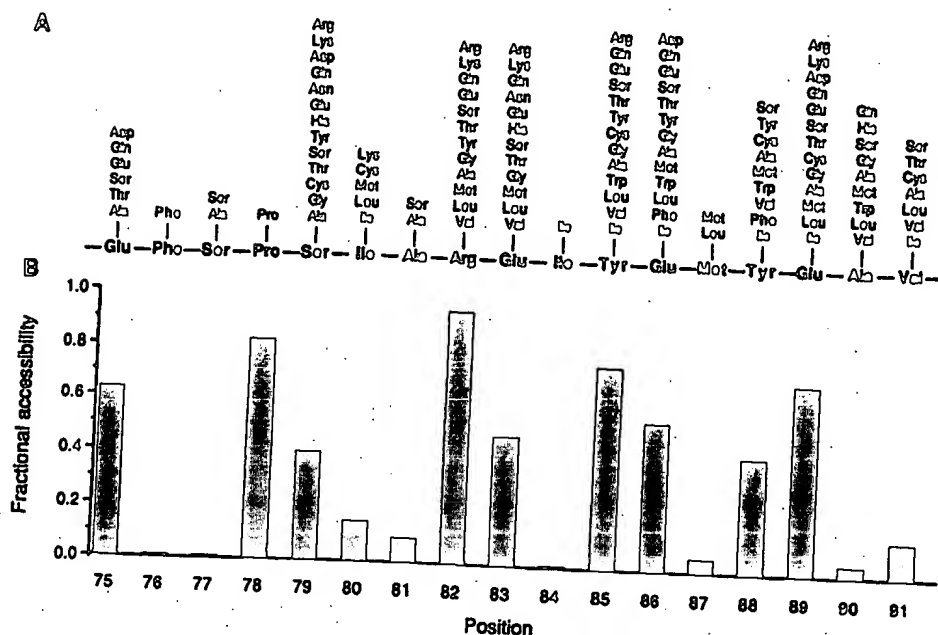
Substitutions at Surface and Buried Positions

In their initial comparisons of the globin sequences, Perutz and co-workers found that most buried residues require nonpolar side chains, whereas few features of surface side chains are generally conserved (6). Similar results have been seen for a number of protein families (2, 4, 5, 7, 17, 18). An example of the sequence tolerance at surface versus buried sites can be seen in Fig. 1, which shows the allowed substitutions in λ repressor at residue positions that are near the dimer interface but distant from the DNA binding surface of the protein (9). These substitutions were identified by a functional

*Authors are in the Department of Biology, Massachusetts Institute of Technology, Cambridge, MA 02139.

Present address: Department of Chemistry and Biochemistry and the Molecular Biology Institute, University of California, Los Angeles, Los Angeles, CA 90024.

Fig. 1. (A) Amino acid substitutions allowed in a short region of λ repressor. The wild-type sequence is shown along the center line. The allowed substitutions shown above each position were identified by randomly mutating one to three codons at a time by using a cassette method and applying a functional selection (9). (B) The fractional solvent accessibility (42) of the wild-type side chain in the protein dimer (43) relative to the same atoms in an Ala-X-Ala model tripeptide.



selection after cassette mutagenesis. A histogram of side chain solvent accessibility in the crystal structure of the dimer is also shown in Fig. 1. At six positions, only the wild-type residue or relatively conservative substitutions are allowed. Five of these positions are buried in the protein. In contrast, most of the highly exposed positions tolerate a wide range of chemically different side chains, including hydrophilic and hydrophobic residues. Hence, it seems that most of the structural information in this region of the protein is carried by the residues that are solvent inaccessible.

Constraints on Core Sequences

Because core residue positions appear to be extremely important for protein folding or stability, we must understand the factors that dictate whether a given core sequence will be acceptable. In general, only hydrophobic or neutral residues are tolerated at buried sites in proteins, undoubtedly because of the large favorable contribution of the hydrophobic effect to protein stability (19). For example, Fig. 2 shows the results of genetic studies used to investigate the substitutions allowed at residue positions that form the hydrophobic core of the NH_2 -terminal domain of λ repressor (20). The acceptable core sequences are composed almost exclusively of Ala, Cys, Thr, Val, Ile, Leu, Met, and Phe. The acceptability of many different residues at each core position presumably reflects the fact that the hydrophobic effect, unlike hydrogen bonding, does not depend on specific residue pairings. Although it is possible to imagine a hypothetical core structure that is stabilized exclusively by residues forming hydrogen bonds and salt bridges, such a core would probably be difficult to construct because hydrogen bonds require pairing of donors and acceptors in an exact geometry. Thus the repertoire of possible structures that use a polar core would probably be extremely limited (21). Polar and charged residues are occasionally found in the cores of proteins, but only at positions where their hydrogen bonding needs can be satisfied (22).

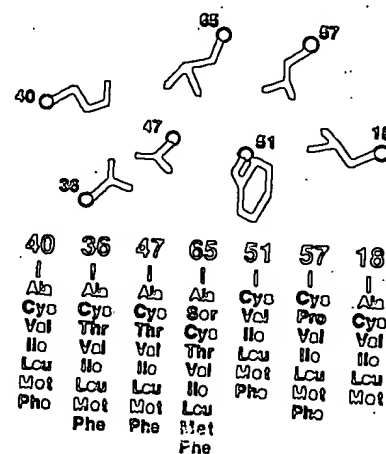
The cores of most proteins are quite closely packed (23), but some volume changes are acceptable. In λ repressor, the overall core volume of acceptable sequences can vary by about 10%. Changes at individual sites, however, can be considerably larger. For example, as shown in Fig. 2, both Phe and Ala are allowed at the same core position in the appropriate sequence contexts. Large volume changes at individual buried sites have also been observed in

phylogenetic studies, where it has been noted that the size decreases and increases at interacting residues are not necessarily related in a simple complementary fashion (5, 7, 17). Rather, local volume changes are accommodated by conformational changes in nearby side chains and by a variety of backbone movements.

The Informational Importance of the Core

With occasional exceptions, the core must remain hydrophobic and maintain a reasonable packing density. However, since the core is composed of side chains that can assume only a limited number of conformations (24), efficient packing must be maintained without steric clashes. How important are hydrophobicity, volume, and steric complementarity in determining whether a given sequence can form an acceptable core? Each factor is essential in a physical sense, as a stable core is probably unable to tolerate unsatisfied hydrogen bonding groups, large holes, or steric overlaps (25). However, in an informational sense, these factors are not equivalent. For example, in experiments in which three core residues of λ repressor were mutated simultaneously, volume was a relatively unimportant informational constraint because three-quarters of all possible combinations of the 20 naturally occurring amino acids had volumes within the range tolerated in the core, and yet most of these sequences were unacceptable (20). In contrast, of the sequences that contained only

Fig. 2. Amino acid substitutions allowed in the core of λ repressor. The wild-type side chains are shown pictorially in the approximate orientation seen in the crystal structure (43). The lists of allowed substitutions at each position are shown below the wild-type side chains. These substitutions were identified by randomly mutating one to four residues at a time by using a cassette method and applying a functional selection (20). Not all substitutions are allowed in every sequence background.



the appropriate hydrophobic residues, a significant fraction were acceptable. Hence, the hydrophobicity of a sequence contains more information about its potential acceptability in the core than does the total side chain volume. Steric compatibility was intermediate between volume and hydrophobicity in informational importance.

The Informational Importance of Surface Sites

We have noted that many surface sites can tolerate a wide variety of side chains, including hydrophilic and hydrophobic residues. This result might be taken to indicate that surface positions contain little structural information. However, Bashford *et al.*, in an extensive analysis of globin sequences (4), found a strong bias against large hydrophobic residues at many surface positions. At one level, this may reflect constraints imposed by protein solubility, because large patches of hydrophobic surface residues would presumably lead to aggregation. At a more fundamental level, protein folding requires a partitioning between surface and buried positions. Consequently, to achieve a unique native state without significant competition from other conformations, it may be important that some sites have a decided preference for exterior rather than interior positions. As a result, many surface sites can accept hydrophobic residues individually, but the surface as a whole can probably tolerate only a moderate number of hydrophobic side chains.

Identification of Residue Roles from Sets of Sequences

Often, a protein of interest is a member of a family of related sequences. What can we infer from the pattern of allowed substitutions at positions in sets of aligned sequences generated by genetic or phylogenetic methods? Residue positions that can accept a number of different side chains, including charged and highly polar residues, are almost certain to be on the protein surface. Residue positions that remain hydrophobic, whether variable or not, are likely to be buried within the structure. In Fig. 3, those residue positions in λ repressor that can accept hydrophilic side chains are shown in orange and those that cannot accept hydrophilic side chains are shown in green. The obligate hydrophobic positions define the core of the structure, whereas positions that can accept hydrophilic side chains define the surface.

Functionally important residues should be conserved in sets of diverse sequences, but it is not possible to decide whether a side chain is functionally or structurally important just because it is invariant or conserved. To make this distinction requires an independent assay of protein folding. The ability of a mutant protein to maintain a stably folded structure can often be measured by biophysical techniques, susceptibility to intracellular proteolysis (26), or by binding to antibodies specific for the native structure (27, 28). In the latter cases, it is possible to screen proteins in mutated clones for the ability to fold even if these proteins are inactive. Sets of sequences that allow formation of a stable structure can then be compared to sets that allow both folding and function, with the active site or binding residues being those that are variable in the set of stable proteins but invariant in the set of functional proteins. The DNA-binding residues of Arc repressor were identified by this method (8). The receptor-binding residues of human growth hormone were also identified by comparing the stabilities and activities of a set of mutant sequences (28). However, in this case, the mutants were generated as hybrid sequences between growth hormone and related hormones with different binding specificities.

Implications for Structure Prediction

At present, the only reliable method for predicting a low-resolution tertiary structure of a new protein is by identifying sequence similarity to a protein whose structure is already known (29, 30). However, it is often difficult to align sequences as the level of sequence similarity decreases, and it is sometimes impossible to detect statistically significant sequence similarity between distantly related proteins. Because the number of known sequences is far greater than the number of known structures, it would be advantageous to increase the reach of the available structural information by improving methods for detecting distant sequence relations and for subsequently aligning these sequences based on structural principles. In a normal homology search, the sequence database is scanned with a single test sequence, and every residue must be weighted equally. However, some residues are more important than others and should be weighted accordingly. Moreover, certain regions of the protein are more likely to contain gaps than others. Both kinds of information can be obtained from sequence sets, and several techniques have

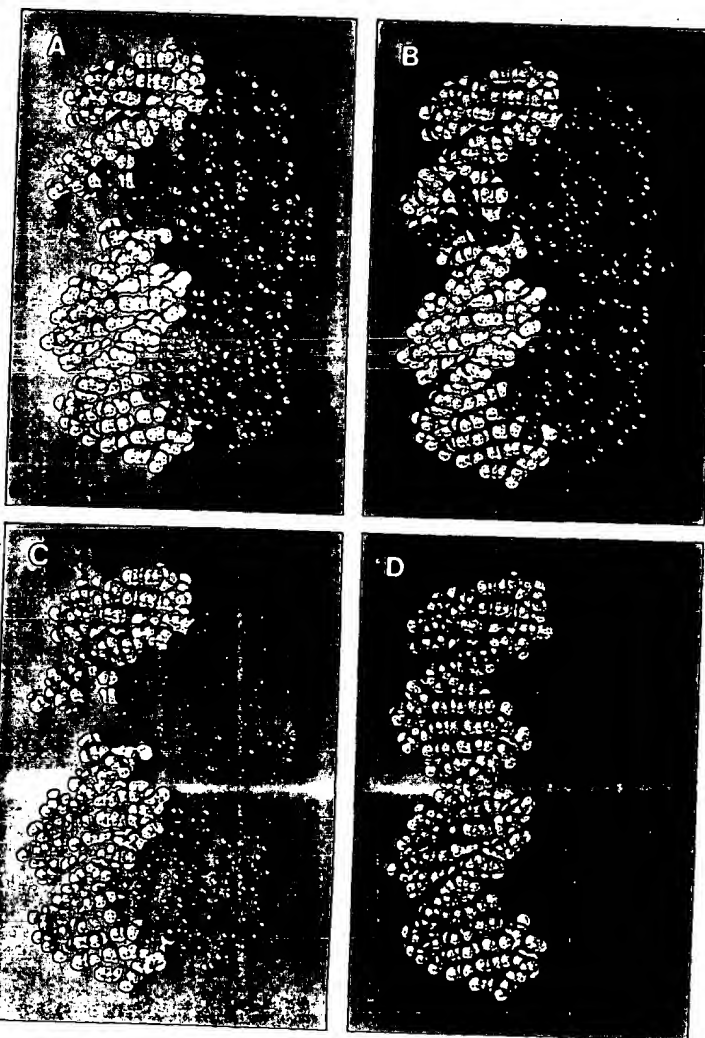


Fig. 3. Tolerance of positions in the NH_2 -terminal domain of λ repressor to hydrophilic side chains. The complex (43) of the repressor dimer (blue) and operator DNA (white) is shown. In (A), positions that can tolerate hydrophilic side chains are shown in orange. The same side chains are shown in (B) without the remaining protein atoms. In (C), positions that require hydrophobic or neutral side chains are shown in green. These side chains are shown in (D) without the remaining protein atoms. About three-fourths of the 92 side chains in the NH_2 -terminal domain are included in both (B) and (D). The remaining positions have not been tested. Data are from (9, 14, 20, 27, 44).

been used to combine such information into more appropriately weighted sequence searches and alignments (31). These methods were used to align the sequences of retroviral proteases with aspartic proteases, which in turn allowed construction of a three-dimensional model for the protease of human immunodeficiency virus type 1 (29). Comparison with the recently determined crystal structure of this protein revealed reasonable agreement in many areas of the predicted structure (32).

The structural information at most surface sites is highly degenerate. Except for functionally important residues, exterior positions seem to be important chiefly in maintaining a reasonably polar surface. The information contained in buried residues is also degenerate, the main requirement being that these residues remain hydrophobic. Thus, at its most basic level, the key structural message in an amino acid sequence may reside in its specific pattern of hydrophobic and hydrophilic residues. This is meant in an informational sense. Clearly, the precise structure and stability of a protein depends on a large number of detailed interactions. It is possible, however, that structural prediction at a more primitive level can be accomplished by concentrating on the most basic informational aspects of an amino acid sequence. For example, amphipathic patterns can be extracted from aligned sets of sequences and used, in some cases, to identify secondary structures.

If a region of secondary structure is packed against the hydrophobic core, a pattern of hydrophobic residues reflecting the periodicity of the secondary structure is expected (33, 34). These patterns can be obscured in individual sequences by hydrophobic residues on the protein surface. It is rare, however, for a surface position to remain hydrophobic over the course of evolution. Consequently, the amphipathic patterns expected for simple secondary structures can be much clearer in a set of related sequences (6). This principle is illustrated in Fig. 4, which shows helical hydrophobic moment plots for the Antennapedia homeodomain sequence (Fig. 4A) and for a composite sequence derived from a set of homologous homeodomain proteins (Fig. 4B) (35). The hydrophobic moment is a simple measure of the degree of amphipathic character of a sequence in a given secondary structure (34). The amphipathic character of the three α -helical regions in the Antennapedia protein (36) is clearly revealed only by the analysis of the combined set of homeodomain sequences. The secondary structure of Arc repressor, a small DNA-binding protein, was recently predicted by a similar method (8) and confirmed by nuclear magnetic resonance studies (37).

The specific pattern of hydrophobic and hydrophilic residues in an amino acid sequence must limit the number of different structures a given sequence can adopt and may indeed define its overall fold. If this is true, then the arrangement of hydrophobic and hydrophilic residues should be a characteristic feature of a particular fold. Sweet and Eisenberg have shown that the correlation of the pattern of hydrophobicity between two protein sequences is a good criterion for their structural relatedness (38). In addition, several studies indicate that patterns of obligatory hydrophobic positions identified from aligned sequences are distinctive features of sequences that adopt the same structure (4, 29, 38, 39). Thus, the order of hydrophobic and hydrophilic residues in a sequence may actually be sufficient information to determine the basic folding pattern of a protein sequence.

Although the pattern of sequence hydrophobicity may be a characteristic feature of a particular fold, it is not yet clear how such patterns could be used for prediction of structure *de novo*. It is important to understand how patterns in sequence space can be related to structures in conformation space. Lau and Dill have approached this problem by studying the properties of simple sequences composed only of H (hydrophobic) and P (polar) groups on two-dimensional lattices (40). An example of such a representa-

tion is shown in Fig. 5. Residues adjacent in the sequence must occupy adjacent squares on the lattice, and two residues cannot occupy the same space. Free energies of particular conformations are evaluated with a single term, an attraction of H groups. By considering chains of ten residues, an exhaustive conformational search for all 1024 possible sequences of H and P residues was possible. For longer sequences only a representative fraction of the allowed sequence or conformation space could be explored. The significant results were as follows: (i) not all sequences can fold into a "native" structure and only a few sequences form a unique native structure; (ii) the probability that a sequence will adopt a unique native structure increases with chain length; and (iii) the native states are compact, contain a hydrophobic core surrounded by polar residues, and contain significant secondary structure. Although the gap between these two-dimensional simulations and three-dimensional structures is large, the use of simple rules and sequence representations yields results similar to those expected for real proteins. Three-dimensional lattice methods are also beginning to be developed and evaluated (41).

Summary

There is more information in a set of related sequences than in a single sequence. A number of practical applications arise from an analysis of the tolerance of residue positions to change. First, such information permits the evaluation of a residue's importance to the function and stability of a protein. This ability to identify the essential elements of a protein sequence may improve our understanding of the determinants of protein folding and stability as well as protein function. Second, patterns of tolerance to amino acid substitutions of varying hydrophilicity can help to identify residues likely to be buried in a protein structure and those likely to occupy

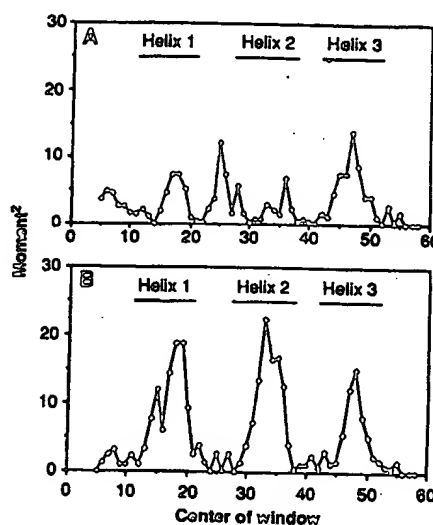


Fig. 4. Helical hydrophobic moments calculated by using (A) the Antennapedia homeodomain sequence or (B) a set of 39 aligned homeodomain sequences (35). The bars indicate the extent of the helical regions identified in nuclear magnetic resonance studies of the Antennapedia homeodomain (36). To determine hydrophobic moments, residues were assigned to one of three groups: H1 (high hydrophobicity = Trp, Ile, Phe, Leu, Met, Val, or Cys); H2 (medium hydrophobicity = Tyr, Pro, Ala, Thr,

His, Gly, or Ser); and H3 (low hydrophobicity = Gln, Asn, Glu, Asp, Lys, or Arg). For the aligned homeodomain sequences, the residues at each position were sorted by their hydrophobicity by using the scale of Fauchere and Pliska (45). Arg and Lys were not counted unless no other residue was found at the position, because they contain long aliphatic side chains and can thereby substitute for nonpolar residues at some buried sites. To account for possible sequence errors and rare exceptions, the most hydrophilic residue allowed at each position was discarded unless it was observed twice. The second most hydrophilic residue was then chosen to represent the hydrophobicity of each position. An eight-residue window was used and the vectors projected radially every 100°. The vector magnitudes were assigned a value of 1, 0, or -1 for positions where the hydrophobicity group was H1, H2, or H3, respectively.

P H P P H P H P H H H P P H

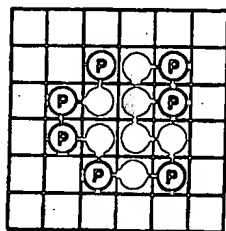


Fig. 5. A representation of one compact conformation for a particular sequence of H and P residues on a two-dimensional square lattice. [Adapted from (40), with permission of the American Chemical Society]

surface positions. The amphipathic patterns that emerge can be used to identify probable regions of secondary structure. Third, incorporating a knowledge of allowed substitutions can improve the ability to detect and align distantly related proteins because the essential residues can be given prominence in the alignment scoring.

As more sequences are determined, it becomes increasingly likely that a protein of interest is a member of a family of related sequences. If this is not the case, it is now possible to use genetic methods to generate lists of allowed amino acid substitutions. Consequently, at least in the short term, it may not be necessary to solve the folding problem for individual protein sequences. Instead, information from sequence sets could be used. Perhaps by simplifying sequence space through the identification of key residues, and by simplifying conformation space as in the lattice methods, it will be possible to develop algorithms to generate a limited number of trial structures. These trial structures could then, in turn, be evaluated by further experiments and more sophisticated energy calculations.

REFERENCES AND NOTES

1. C. J. Epstein, R. F. Goldberger, C. B. Anfinsen, *Cold Spring Harbor Symp. Quant. Biol.* 28, 439 (1963); C. B. Anfinsen, *Science* 181, 223 (1973).
2. R. E. Dickerson, *Sci. Am.* 242, 136 (March 1980).
3. M. D. Hampsey, G. Das, F. Sherman, *FEBS Lett.* 231, 275 (1988).
4. D. Bashford, C. Chothia, A. M. Lesk, *J. Mol. Biol.* 196, 199 (1987).
5. A. M. Lesk and C. Chothia, *ibid.* 136, 225 (1980).
6. M. F. Perutz, J. C. Kendrew, H. C. Watson, *ibid.* 13, 669 (1965).
7. C. Chothia and A. M. Lesk, *Cold Spring Harbor Symp. Quant. Biol.* 52, 399 (1965).
8. J. U. Bowie and R. T. Sauer, *Proc. Natl. Acad. Sci. U.S.A.* 86, 2152 (1989).
9. J. F. Reidhaar-Olson and R. T. Sauer, *Science* 241, 53 (1988); *Protein Struct. Funct. Genet.*, in press.
10. D. Shortle, *J. Biol. Chem.* 264, 5315 (1989).
11. J. H. Müller et al., *J. Mol. Biol.* 131, 191 (1979).

12. S. Sprang et al., *Science* 237, 905 (1987); C. S. Craik, S. Rocznick, C. Largman, W. J. Rutter, *ibid.*, p. 909.
13. H. C. M. Nelson and R. T. Sauer, *J. Mol. Biol.* 192, 27 (1986).
14. M. H. Hecht, J. M. Sturtevant, R. T. Sauer, *Proc. Natl. Acad. Sci. U.S.A.* 81, 5685 (1984).
15. T. Alber, D. Sun, J. A. Nye, D. C. Muchmore, B. W. Matthews, *Biochemistry* 26, 3754 (1987).
16. D. Shortle and A. K. Meeker, *Protein Struct. Funct. Genet.* 1, 81 (1986).
17. A. M. Lesk and C. Chothia, *J. Mol. Biol.* 160, 325 (1982).
18. W. R. Taylor, *ibid.* 188, 233 (1986).
19. W. Kauzmann, *Adv. Protein Chem.* 14, 1 (1959); R. L. Baldwin, *Proc. Natl. Acad. Sci. U.S.A.* 83, 8069 (1986).
20. W. A. Lim and R. T. Sauer, *Nature* 339, 31 (1989); in preparation.
21. Lesk and Chothia (5) have argued that a protein core composed solely of hydrogen-bonded residues would also be inviable on evolutionary grounds, as a mutational change in one core residue would require compensating changes in any interacting residue or residues to maintain a stable structure.
22. T. M. Gray and B. W. Matthews, *J. Mol. Biol.* 175, 75 (1984); E. N. Baker and R. E. Hubbard, *Prog. Biophys. Mol. Biol.* 44, 97 (1984).
23. F. M. Richards, *J. Mol. Biol.* 92, 1 (1974).
24. J. W. Ponder and F. M. Richards, *ibid.* 193, 775 (1987).
25. J. T. Kellis, Jr., K. Nyberg, A. R. Fersht, *Biochemistry* 28, 4914 (1989); W. S. Sandberg and T. C. Terwilliger, *Science* 245, 54 (1989).
26. A. A. Pakula and R. T. Sauer, *Protein Struct. Funct. Genet.* 5, 202 (1989).
27. B. C. Cunningham and J. A. Wells, *Science* 244, 1081 (1989); R. M. Breyer and R. T. Sauer, *J. Biol. Chem.* 264, 13348 (1989).
28. B. C. Cunningham, P. Jhurani, P. Ng, J. A. Wells, *Science* 243, 1330 (1989).
29. L. H. Pearl and W. R. Taylor, *Nature* 329, 351 (1987).
30. W. J. Brown et al., *J. Mol. Biol.* 42, 65 (1969); J. Groer, *ibid.* 153, 1027 (1981); J. M. Berg, *Proc. Natl. Acad. Sci. U.S.A.* 85, 99 (1988).
31. W. R. Taylor, *Protein Eng.* 2, 77 (1988).
32. M. A. Navia et al., *Nature* 337, 615 (1989).
33. M. Schiffer and A. B. Edmundson, *Biophys. J.* 7, 121 (1967); V. I. Lim, *J. Mol. Biol.* 88, 857 (1974); *ibid.*, p. 873.
34. D. Eisenberg, R. M. Weiss, T. C. Terwilliger, *Nature* 299, 371 (1982); D. Eisenberg, D. Schwarz, M. Komaromy, R. Wall, *J. Mol. Biol.* 179, 125 (1984); D. Eisenberg, R. M. Weiss, T. C. Terwilliger, *Proc. Natl. Acad. Sci. U.S.A.* 81, 140 (1984).
35. T. R. Burglin, *Cell* 53, 339 (1988).
36. G. Otting et al., *EMBO J.* 7, 4305 (1988).
37. J. N. Breg, R. Boelens, A. V. E. George, R. Kaptein, *Biochemistry* 28, 9826 (1989); M. G. Zagorski, J. U. Bowie, A. K. Vershon, R. T. Sauer, D. J. Parel, *ibid.*, p. 9813.
38. R. M. Sweet and D. Eisenberg, *J. Mol. Biol.* 171, 479 (1983).
39. J. U. Bowie, N. D. Clarke, C. O. Pabo, R. T. Sauer, *Protein Struct. Funct. Genet.*, in preparation.
40. K. F. Lau and K. A. Dill, *Macromolecules* 22, 3986 (1989).
41. A. Sikorski and J. Skolnick, *Proc. Natl. Acad. Sci. U.S.A.* 86, 2668 (1989); A. Kolinski, J. Skolnick, R. Yaris, *Biopolymers* 26, 937 (1987); D. G. Cowell and R. L. Jernigan, *Biochemistry*, in press.
42. B. Lee and F. M. Richards, *J. Mol. Biol.* 55, 379 (1971).
43. S. R. Jordan and C. O. Pabo, *Science* 242, 893 (1988).
44. R. M. Breyer, thesis, Massachusetts Institute of Technology, Cambridge (1988).
45. J.-L. Fauchere and V. Pliska, *Eur. J. Med. Chem.-Chim. Ther.* 18, 369 (1983).
46. We thank C. O. Pabo and S. Jordan for coordinates of the NH₂-terminal domain of λ repressor and its operator complex. We also thank P. Schimmel for the use of his graphics system and J. Burnbaum and C. Franclyn for assistance. Supported in part by NIH grant AI-15706 and predoctoral grants from NSF (J.R.-O.) and Howard Hughes Medical Institute (W.A.L.).

**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☒ **BLACK BORDERS**
- ☐ **IMAGE CUT OFF AT TOP, BOTTOM OR SIDES**
- ☐ **FADED TEXT OR DRAWING**
- ☐ **BLURRED OR ILLEGIBLE TEXT OR DRAWING**
- ☐ **SKEWED/SLANTED IMAGES**
- ☐ **COLOR OR BLACK AND WHITE PHOTOGRAPHS**
- ☐ **GRAY SCALE DOCUMENTS**
- ☐ **LINES OR MARKS ON ORIGINAL DOCUMENT**
- ☐ **REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY**
- ☐ **OTHER:** _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.